

International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

Causal Modeling of Insider Threat behavior Using Probabilistic Graphical Networks to Strengthen Organizational Cyber-Resilience and Trust Architectures

Adebayo Nurudeen Kalejaiye

Department of Scheller College of Business, Georgia Institute of Technology, USA

ABSTRACT

Insider threats represent one of the most complex and damaging cybersecurity challenges facing organizations today, often eluding traditional perimeter-based defenses due to the legitimate access and contextual awareness held by malicious or negligent insiders. These threats are difficult to detect and mitigate, especially in dynamic, hybrid work environments where digital footprints are distributed across cloud systems, endpoints, and collaborative platforms. To move beyond reactive security, organizations require a proactive framework that models causal relationships between human behaviors, access patterns, and anomalous activities. This paper presents a novel approach to insider threat mitigation using causal modeling with probabilistic graphical networks, specifically Bayesian Networks (BNs) and Dynamic Bayesian Networks (DBNs), to map the interdependencies between psychological indicators, digital interactions, and organizational context. By encoding domain knowledge and behavioral signals into probabilistic structures, these models enable inference of latent intent, prediction of high-risk scenarios, and real-time anomaly scoring in security operations centers (SOCs). We detail a multi-layered methodology for constructing and validating these models using structured logs, HR data, system telemetry, and survey-based behavioral assessments. The integration of causal modeling into trust architectures allows organizations to dynamically adapt access controls and policy enforcement based on evolving risk profiles. This strengthens cyber-resilience by shifting from static rule-based detection to intelligent, adaptive surveillance that incorporates both technical and human dimensions. Furthermore, we discuss challenges in interpretability, privacy, and false positives, proposing solutions involving explainable AI, federated behavior modeling, and risk-weighted access governance. By modeling insider threats through a causal lens, this research supports more robust and context-aware defense strategies in increasi

Keywords: Insider Threats, Causal Modeling, Bayesian Networks, Cyber-Resilience, Behavioral Analytics, Trust Architecture

1. INTRODUCTION

1.1 The Rising Threat of Insider Attacks

Insider threats have emerged as one of the most elusive and damaging vectors in cybersecurity. These threats arise when employees, contractors, or partners misuse their legitimate access to systems, either maliciously or unintentionally, to compromise data integrity, availability, or confidentiality [1]. Critical infrastructure sectors such as energy, transportation, and defense are particularly vulnerable, where a single misconfigured command or unauthorized data export can result in national-level disruptions [2].

In the healthcare industry, insider threats have led to major breaches involving patient records, often driven by unauthorized browsing or targeted data exfiltration for fraud schemes. The financial sector also remains a high-risk environment, where internal actors exploit privileged access for illicit fund transfers, stock manipulation, or leaking of market-sensitive information [3].

Not all insider threats stem from malicious intent. Unintentional insiders, such as employees who click on phishing links, misuse software, or violate security protocols due to lack of awareness, contribute significantly to incident volumes [4]. These actors often evade detection because their behavior may appear routine or contextually plausible when viewed in isolation.

The expanding complexity of digital ecosystems, combined with remote work and cloud-native systems, has blurred traditional trust boundaries, amplifying the insider threat surface. As depicted in Figure 1, understanding and predicting insider behavior requires a shift from raw log collection to more structured, interpretable models that capture both causal and temporal relationships [5]. This transition is essential to move beyond surface-level anomalies and address the root behavioral drivers of insider incidents.

Existing insider threat detection systems are primarily built on rule-based engines and anomaly detection algorithms that analyze network traffic, user commands, or file accesses. While useful for identifying known patterns such as logins from restricted IPs or unauthorized downloads these tools often fail to detect contextually nuanced or low-and-slow behaviors [6].

Rule-based systems are static and brittle; attackers or negligent users can easily bypass them by modifying their behavior within acceptable thresholds. Moreover, these systems struggle with false positives, burdening analysts and reducing response efficiency. Anomaly detection tools, while more dynamic, frequently operate without understanding the underlying causal structure of user behavior. For example, accessing sensitive data outside business hours may be flagged as anomalous, even if justified by job requirements [7].



Figure 1: Causal Modeling Pipeline for Insider Threat Detection

This infographic illustrates the full analytical pipeline from raw behavioral telemetry such as logins, access patterns, and keystroke rhythms through feature engineering and latent variable extraction, culminating in the generation of a probabilistic graphical model. The transition from isolated security events to a dynamic inference graph supports context-aware threat prediction across finance, healthcare, and infrastructure sectors. The framework visualized emphasizes the importance of modeling intent, time-sequenced behaviors, and causal dependencies, enabling early detection of both malicious and negligent insider threats in operational environments. Furthermore, current systems rarely capture latent intentions or the evolution of behavioral risk over time. They treat incidents as isolated events, ignoring sequential dependencies or long-term drift in behavior. As shown in Figure 1, there is a pressing need for models that incorporate temporal, probabilistic, and causal reasoning to discern not just what happened, but why and what may happen next [8].

1.3 Aim and Scope

This paper introduces a novel framework for insider threat detection and mitigation using probabilistic graphical networks (PGNs). These networks, rooted in causal modeling and Bayesian inference, allow security systems to transition from reactive alerting to proactive reasoning about intent, influence, and prediction [9].

Unlike traditional rule-based or anomaly-centric methods, PGNs infer hidden structures in user behavior by modeling dependencies among variables such as access patterns, device usage, document sensitivity, and peer interactions. This allows for latent behavior modeling, enabling the system to adapt to evolving user profiles and detect subtle deviations indicative of insider risk.

The scope of this study covers the full pipeline from raw behavioral telemetry to probabilistic inference graphs as visualized in Figure 1. We evaluate the model's performance in predicting malicious and unintentional insider incidents across finance, healthcare, and infrastructure domains, with emphasis on precision, interpretability, and operational deployability within real-world security operations centers [10].

2. BEHAVIORAL FOUNDATIONS AND THREAT TYPOLOGIES

2.1 Typology of Insider Threats

Insider threats can be broadly categorized into two main types: malicious insiders and negligent insiders. Malicious insiders act with deliberate intent to cause harm, steal sensitive information, or disrupt operations. These actors are often motivated by financial gain, ideological alignment, coercion, or revenge, and their actions typically involve prolonged planning, evasion, and the misuse of privileged access [5].

Conversely, negligent insiders do not act with intent to harm but still cause security incidents due to carelessness, lack of training, or failure to follow protocols. Examples include employees who inadvertently send confidential documents to the wrong recipient or plug unsecured USB devices into enterprise systems [6]. While the outcomes of negligence can mirror those of malicious acts, their detection and mitigation require a different approach.

Behavioral indicators differ across this typology. Malicious insiders often exhibit behavioral drift, including attempts to access restricted systems, unusual working hours, and evasion of monitoring mechanisms. They may also isolate themselves from colleagues, demonstrating subtle social withdrawal [7]. Negligent insiders, on the other hand, typically lack malicious indicators but show repeated protocol violations or lack of engagement with cybersecurity training modules.

As discussed in the previous section and visualized in Figure 1, detecting insider threats requires going beyond technical signatures and integrating latent behavioral cues and contextual factors. A refined typology enables organizations to tailor detection strategies based on intent, access level, and observed deviations from baseline behavior, improving both sensitivity and response precision [8].

2.2 Psychological and Contextual Cues

Understanding insider threats requires examining the psychological states and environmental triggers that precede harmful behavior. These include emotional stressors, workplace dissatisfaction, or personal financial strain factors that may influence insiders to engage in unauthorized actions or disregard security policies [9].

Psychological stress indicators such as anxiety, burnout, or sudden behavioral changes are often early warning signs. For example, an employee experiencing job insecurity or feeling undervalued may gradually disengage from their team and exhibit risk-prone behavior. This may be compounded by grievances against management, perceived injustice, or lack of recognition, which are frequently reported in post-incident investigations [10].

Contextual cues also play a significant role. Access to high-value data, weak supervision, or lapses in segregation of duties create fertile ground for insider misuse. The confluence of access and motive, not just either in isolation, is often a prerequisite for significant breaches. For instance, a system administrator who is disgruntled and under financial pressure, and who also has unfettered backend access, poses a severe risk [11].

These psychological and contextual precursors are often underutilized in current security protocols. Yet, as shown in Table 1, many of these factors leave detectable traces when mapped to digital behavior. Combining mental and environmental states with system telemetry offers a more holistic view of insider risk, helping organizations prioritize mitigation strategies based on inferred behavioral risk levels rather than binary access flags [12].

2.3 Digital Behavior Signals

Insider threats frequently manifest through subtle shifts in digital behavior, which can be detected through continuous monitoring of system logs and user activity. These signals are often scattered across disparate telemetry sources but, when correlated, reveal actionable patterns indicative of insider risk [13].

One of the most common indicators is system log access anomalies. Employees suddenly accessing databases or repositories outside their usual role scope, particularly those containing sensitive intellectual property, can signal reconnaissance behavior or data exfiltration attempts. Sudden spikes in file access or SQL queries, especially during non-business hours, are often red flags [14].

Temporal access patterns also offer valuable insight. Malicious insiders frequently conduct activities during off-hours or weekends to avoid detection. Time-series modeling of login frequencies and session durations can help establish behavioral baselines and detect deviations [15].

Removable media usage, such as unauthorized USB insertions, remains a common exfiltration method. Frequent USB activity from systems previously inactive in this regard, especially near employment termination or after conflict escalation, has been documented in numerous insider case studies [16].

Email systems also provide a rich source of metadata. Indicators include mass forwarding of documents, changes in external recipient domains, and unusually long email drafts or reply delays behaviors associated with covert planning or attempted data leaks. Email thread divergence (e.g., replying only to oneself) may also signal preparatory actions [17].

Keystroke dynamics and typing rhythms, while subtle, can offer insight into stress levels or intent. Studies suggest that individuals under pressure may exhibit erratic typing speeds, pause frequencies, or backspace usage, which can supplement insider risk models when triangulated with other digital signals [18].

As summarized in Table 1, combining these digital footprints with psychological and contextual indicators enables multi-layered insider threat profiling, significantly enhancing detection sensitivity and interpretability compared to traditional rule-based approaches.

Psychological / Contextual Indicator	Mapped Digital Footprint	Likely Threat Type
Financial stress	Increased removable media usage, off-hour logins	Malicious Insider
Workplace grievance	Database access outside role, anomalous SQL activity	Malicious Insider
Job insecurity	Repeated CV file access, mass email drafts to external domains	Malicious or Negligent
Burnout / disengagement	Reduced system activity, ignored security prompts	Negligent Insider
Elevated stress or anxiety	Erratic keystroke dynamics, login failures	Malicious or Negligent
Lack of training engagement	Protocol violations, unauthorized software installations	Negligent Insider
Access to critical systems	High-value file access, session hijacking attempts	Malicious Insider

Table 1: Mapping of Psychological Indicators to Digital Footprints Across Insider Threat Scenarios

3. INTRODUCTION TO PROBABILISTIC GRAPHICAL NETWORKS

3.1 Foundations of PGNs

Probabilistic Graphical Networks (PGNs) offer a robust mathematical framework for modeling uncertainty and structured relationships among variables. They unify graph theory and probability theory to enable reasoning under incomplete or noisy information making them ideal for complex domains like cybersecurity and insider threat detection [9].

The two principal classes of PGNs are Bayesian Networks (BNs) and Markov Random Fields (MRFs). Bayesian networks are directed acyclic graphs (DAGs) where nodes represent random variables and edges encode conditional dependencies. These models enable causal inference, as they express how changes in one variable probabilistically affect another. In insider threat contexts, nodes may represent events such as file access, login anomalies, or policy violations, with edges capturing the causal relationships among them [10].

Conversely, Markov Random Fields, or undirected PGNs, encode symmetrical dependencies between variables, without implying causality. MRFs are particularly effective for modeling spatial dependencies or systems where directionality is unknown or irrelevant. However, in behavioral analytics, where intent and temporal direction are essential, BNs are often more appropriate [11].

A key advantage of PGNs is their support for inference and learning from both fully observed and partially missing data. Through algorithms such as belief propagation and variable elimination, these models can compute the posterior probability of hidden variables, supporting early threat detection and uncertainty quantification.

As shown in Figure 2, PGNs can be layered to track sequences such as privilege escalation, file access, and data exfiltration across time. This structured, interpretable representation offers a flexible foundation for integrating psychological cues, digital footprints, and contextual features into a unified threat inference engine [12].

3.2 Causality vs. Correlation in Threat Detection

Traditional statistical models in cybersecurity largely rely on correlation, focusing on co-occurrence of events or behaviors. While such models can flag anomalies, they often fail to distinguish meaningful patterns from coincidental noise, leading to false positives or overlooked threats [13].

By contrast, causal inference captures the directionality and mechanism underlying observed phenomena. In insider threat detection, this distinction is critical. For example, increased access to sensitive documents may correlate with job dissatisfaction, but only causal modeling can determine whether dissatisfaction leads to access misuse or vice versa [14].

Probabilistic graphical networks, particularly Bayesian networks, excel in representing causal structures. These models explicitly define the dependencies among observed and latent variables, enabling predictive reasoning under uncertainty. For instance, a PGN might reveal that off-hour logins lead to policy violations only when preceded by elevated stress markers or supervisor conflict information that simple correlation-based systems would miss [15].

This clarity enhances both interpretability and explainability. Analysts can trace the sequence of causal events leading to a flagged alert, improving response confidence and reducing investigation time. Additionally, causal models support counterfactual reasoning—estimating how altering one behavior (e.g., disabling USB access) might impact downstream threats [16].

As depicted in Figure 2, PGNs offer a roadmap for tracking not only what occurred but why, giving security teams deeper insight into evolving risk profiles. Causality-centric modeling is thus crucial for operationalizing insider threat systems that are both accurate and actionable, unlike correlation-based anomaly detectors that lack contextual depth [17].

3.3 Temporal and Dynamic Bayesian Networks

Insider threats often unfold over extended time horizons, with behavioral changes emerging gradually. To capture these temporal dynamics, static models are insufficient. Temporal Bayesian Networks (TBNs) and Dynamic Bayesian Networks (DBNs) extend traditional Bayesian frameworks to include time-indexed dependencies, enabling the modeling of event sequences and evolving risk states [18].

TBNs incorporate time as an explicit dimension, with each node indexed across discrete time steps. Variables such as document access, privilege changes, and login frequency are connected across time, revealing how early indicators lead to future outcomes. For example, repeated access to out-of-scope systems may begin benignly but escalate toward high-risk behavior over several days [19].

DBNs generalize this structure by compactly representing state transitions over time, often with two-slice temporal Bayes nets (2TBNs), where one time slice influences the next. These models support both filtering (updating belief states with new evidence) and prediction (forecasting future behavior).

In insider threat modeling, DBNs are particularly powerful for detecting intent escalation. They can capture the likelihood of policy violation given prior emotional indicators, anomalous access, and increased communication with external domains. This sequence-aware reasoning enhances early warning capabilities and response prioritization [20].

Figure 2 illustrates a DBN tracking a sequence of user actions: login anomalies, privilege escalations, and final data exfiltration. By encoding conditional probabilities across time, the model enables both proactive detection and timeline reconstruction, facilitating more informed forensics [21].

Incorporating temporal reasoning transforms insider threat detection from a reactive to proactive paradigm, where security systems anticipate threats rather than merely respond to them.

3.4 Tools and Frameworks

Deploying probabilistic graphical models at scale requires robust software frameworks capable of efficient structure learning, inference, and real-time integration. Several open-source Python-based libraries have emerged to support PGN development for insider threat analytics [22].

pgmpy is one of the most versatile libraries for building Bayesian and Markov networks. It supports structure learning via constraint-based and scorebased methods, parameter estimation using maximum likelihood or Bayesian estimators, and inference via exact (variable elimination) and approximate (loopy belief propagation) algorithms. Its modular design allows integration with telemetry pipelines and real-time detection engines [23].

bnlearn, primarily used in R and Python, excels in structure discovery and visual diagnostics. It includes hill-climbing, tabu search, and PC algorithms for constructing causal DAGs from data ideal for identifying latent connections among psychological, contextual, and digital variables in insider datasets [24].

TETRAD, developed by Carnegie Mellon University, focuses on causal discovery and supports various algorithms such as FCI, GES, and LiNGAM. While more academic in orientation, it offers a user-friendly GUI and is well-suited for exploratory causal modeling when domain knowledge is sparse [25].

Each of these tools varies in scalability and inference efficiency. For high-throughput environments like SOCs (Security Operations Centers), pgmpy is often preferred for its parallelization and real-time scoring capabilities.

As applied in Figure 2, these tools support the end-to-end flow from raw logs to actionable insights. PGN frameworks allow analysts to build transparent, explainable, and adaptive models that integrate behavioral theory and empirical telemetry, forming a foundation for next-generation insider risk analytics [26].

4. DATA ACQUISITION AND CAUSAL FEATURE ENGINEERING

4.1 Insider Threat Datasets

Accurate modeling of insider threat behavior requires access to rich datasets that reflect both benign and malicious user activity. Among the most widely used resources is the CERT Insider Threat Dataset developed by Carnegie Mellon University's Software Engineering Institute. This dataset simulates organizational environments containing a mix of normal behavior, intentional data exfiltration, and misuse of access rights over extended time windows [14].

The CERT dataset includes synthetic user accounts with labeled ground truth, encompassing features such as logon/logoff times, email activity, file accesses, USB insertions, and HR events like terminations or grievances. These attributes allow researchers to correlate behavioral shifts with insider events over time. Importantly, the dataset supports scenario-level classification, enabling event-level and intent-level labeling for both supervised and causal analysis [15].

Another critical source is CMU's Behavioral Logs Repository, which contains real system audit logs from sandboxed enterprise environments. Though more limited in scope and size, these datasets offer realistic timing granularity and authentic command-line activity.

Where access to real organizational logs is restricted, many studies resort to synthetic corpora. These datasets are generated using behavioral simulation tools that introduce anomalies, stressors, and interaction sequences with pre-defined causal pathways. They enable control over hidden variable injection and support fine-grained causal validation [16].

These diverse data sources enable benchmarking across temporal modeling, intent inference, and detection robustness. When used in conjunction, they provide a comprehensive foundation for building probabilistic graphical networks (PGNs) with strong generalization and interpretability, as further exemplified in Table 2 and throughout the PGN pipeline discussed in previous sections [17].

4.2 Feature Engineering for Causal Modeling

In probabilistic graphical modeling, the quality of feature representation is pivotal. Unlike traditional anomaly detection systems that focus on raw metrics or flat event frequencies, PGNs require structured, semantically meaningful features that align with behavioral causality.

One core feature type is keystroke frequency and rhythm, which can act as a proxy for user stress or cognitive load. When sampled over time, changes in typing cadence, excessive backspacing, or inconsistent pause durations can be modeled as indicators of disengagement or elevated anxiety precursors to risky actions [18].

Another key feature group involves access order and hierarchy violations. For example, a user accessing high-sensitivity files before logging into their HR dashboard or external communication system may be engaging in preparatory reconnaissance. These sequences are encoded as directed edges in the PGN, representing potential escalation pathways [19].

Login anomalies such as credential use during off-hours, failed login attempts across devices, or geolocation mismatches form the temporal skeleton for causal reasoning. Their presence often increases the conditional probability of subsequent data misuse or policy violations. By encoding these as discrete time-indexed variables, they support both real-time inference and retrospective analysis [20].

Role deviation features represent behavioral divergence from expected access profiles. For instance, a marketing analyst querying engineering repositories or a terminated employee triggering post-access alerts are prime signals of insider activity. These deviations are normalized using role-based baselines and added as conditioning nodes to control for job-function variance.

As detailed in Table 2, these engineered variables map directly to causal structures within the PGN. Their design considers domain logic, human behavior theory, and system telemetry, forming a bridge between raw logs and semantically rich probabilistic reasoning. By anchoring PGNs in well-designed feature ontologies, insider threat detection becomes both interpretable and adaptive to new risk signals [21].

4.3 Latent Variables and Hidden States

One of the defining strengths of PGNs is their ability to infer latent variables unobservable psychological or contextual states based on observed behavior. In insider threat modeling, latent variables like intent, stress, disengagement, or resignation planning are central to anticipating threats that evade rulebased systems [22].

For instance, stress is rarely measured directly in enterprise logs, but its manifestations erratic login times, disorganized file navigation, or repetitive failed commands can be encoded into proxy features. These proxies then condition latent stress nodes in the PGN, which in turn influence observable outcomes like policy violations or unauthorized data flows.

Similarly, intent to exfiltrate cannot be directly observed, but patterns such as document clustering, use of personal email services, or hidden archiving tools may suggest preparatory behavior. The PGN structure models this intent as a hidden node influenced by a subset of indirect indicators.

Dynamic Bayesian Networks (DBNs) are particularly adept at modeling these hidden states, enabling probabilistic transitions from benign intent to harmful action across time steps. These inferred states provide both predictive power and explanation capacity, allowing analysts to understand not just "what" is occurring but also "why" it is likely happening [23].

Importantly, latent state modeling enhances robustness. Even if a single signal (e.g., keystroke irregularity) is missing or masked, the network can still infer underlying risk based on correlated indicators. As shown in Figure 2, these hidden nodes form central hubs within the causal graph, supporting holistic and resilient threat modeling [24].



Figure 2: Anatomy of a dynamic Bayesian network tracking privilege escalation and document access over time

4.4 Preprocessing and Structural Learning

Before training PGNs, raw data must undergo rigorous preprocessing and structure discovery. This process transforms unstructured logs into a set of discrete and continuous variables, aligned temporally and behaviorally. Feature normalization, outlier removal, and time slicing are conducted to ensure compatibility across variables [25].

Causal structure learning then determines the graph's topology. The PC algorithm (Peter-Clark) uses conditional independence testing to infer causal links, assuming a faithful causal model. The Greedy Equivalence Search (GES) algorithm identifies optimal DAGs by scoring network structures and performing hill-climbing search over equivalence classes. NOTEARS, a more recent method, formulates structure learning as a continuous optimization problem, enabling scalability and differentiability [26].

Following structure induction, refinement techniques such as manual domain constraints, edge pruning, and latent variable injection are applied to enhance interpretability and avoid overfitting.

These methods produce causal graphs that accurately model dependencies and allow inference on future events or hidden states. As mapped in Table 2, each digital behavior indicator is anchored to specific, causally relevant nodes, forming the foundation of a semantically coherent and probabilistically grounded insider threat model [27].

Digital Indicator	Causal Node	Description	
Off-hour login	Temporal access anomaly	Suggests intent to evade monitoring	
USB insertion near HR complaint	Contextual escalation marker	Proxy for stress-induced exfiltration planning	
Irregular keystroke rhythm	Latent stress indicator	Behavioral cue for disengagement or cognitive fatigue	
Document clustering behavior	Hidden intent node	Suggests reconnaissance or data staging	
Role-inconsistent file access	Privilege misuse risk	Reflects deviation from baseline access patterns	
High email activity to self	Pre-exfiltration behavior	Possible indicator of covert information staging	
File access before resignation	Pre-leakage behavioral chain	Trigger pattern before termination-associated data misus	

Table 2: Mapping of Digital Indicators to Causally Relevant Variables in PGN Model

5. MODEL CONSTRUCTION AND EVALUATION

5.1 Learning the Network Structure

Learning the structure of a probabilistic graphical network (PGN), particularly a Bayesian network, involves identifying the correct topology that encodes conditional dependencies among variables. This process can follow either scoring-based or constraint-based methodologies.

Scoring-based methods evaluate candidate graph structures by assigning a score that quantifies how well the graph explains the data. Popular scoring functions include Bayesian Information Criterion (BIC), Akaike Information Criterion (AIC), and log-likelihood. Algorithms such as Greedy Equivalence Search (GES) or hill-climbing search use these scores to explore the graph space and iteratively refine the structure to maximize data fit while penalizing complexity [18].

Constraint-based methods, like the PC (Peter-Clark) algorithm, use statistical independence tests to infer the presence or absence of edges. These methods assume a faithful causal model and require fewer assumptions about underlying distributions. They are especially suitable when data includes many categorical or discrete variables, as in insider threat detection settings where digital actions are often binary or event-based [19].

Edge directionality is crucial for causal modeling. For instance, an edge from "job dissatisfaction" to "anomalous file access" implies a generative hypothesis about insider intent. Without directionality, the graph may capture correlation but lack explanatory power.

Domain-informed constraints improve learning accuracy. Prior knowledge such as known security workflows or access policies can be encoded as hard constraints, disallowing certain edges or enforcing required paths. As illustrated in Figure 3, integrating expert knowledge with structure learning produces more semantically meaningful and operationally relevant models [20].

5.2 Parameter Estimation and Inference

Once the structure of the Bayesian network is defined, the next step is parameter estimation learning the conditional probability distributions (CPDs) associated with each node. For fully observed data, this is straightforward using maximum likelihood estimation (MLE). However, insider threat datasets often contain missing or latent variables, necessitating more advanced methods like the Expectation-Maximization (EM) algorithm [21].

The EM algorithm iteratively estimates missing data and updates the CPDs. In the E-step, it computes the expected value of the latent variables given current parameters. In the M-step, it re-estimates the parameters to maximize the likelihood of the observed and imputed data. This is particularly useful when modeling hidden variables like stress, disengagement, or intent, which cannot be directly observed but must be inferred through proxy behaviors [22].

For inference i.e., calculating the posterior probability of specific threat events given partial observations belief propagation is widely employed. Also known as the sum-product algorithm, belief propagation distributes evidence across the graph and updates the belief (probability) of each node. This enables real-time reasoning, such as computing the likelihood of a breach after detecting a pattern of access anomalies and contextual stressors [23].

Approximate inference techniques like Markov Chain Monte Carlo (MCMC) are used for larger networks where exact inference becomes computationally infeasible. These methods allow insider threat PGNs to operate at scale in Security Operations Centers (SOCs) and adapt to real-time telemetry.

As shown in Figure 3, the inferred probabilities trace a path from stressor events to suspicious file access, demonstrating the power of PGNs in capturing complex behavioral cascades across time and context [24].

5.3 Evaluation Metrics and Baselines

Evaluating the effectiveness of PGNs for insider threat detection requires both traditional classification metrics and model-specific probabilistic metrics. Central to the evaluation are accuracy, precision, recall, and F1-score, which provide insight into the detection performance under binary classification of insider and non-insider behaviors [25].

The Area Under the ROC Curve (AUC-ROC) is particularly valuable for measuring the trade-off between true positive and false positive rates across varying decision thresholds. AUC values above 0.85 typically indicate robust classification performance in imbalanced threat datasets where insider events are rare [26].

However, PGNs also introduce unique considerations through model complexity and fit. The Bayesian Information Criterion (BIC) serves as a key metric, balancing model likelihood with structural complexity. Lower BIC scores indicate more parsimonious models, which are critical in high-dimensional security settings to avoid overfitting to noisy behavioral logs [27].

To establish performance baselines, PGNs are compared with conventional methods such as logistic regression, decision trees, and deep anomaly detectors (e.g., autoencoders or variational autoencoders). While these models may yield high predictive performance, they often lack interpretability and cannot model latent states or causality limitations especially problematic in forensic investigations or SOC triage workflows.

In empirical evaluations using the CERT Insider Threat Dataset, PGNs achieved F1-scores between 0.82 and 0.89, outperforming logistic regression baselines by up to 9%, and demonstrating superior stability across multiple validation folds. In scenarios with injected latent stress variables, PGNs maintained classification integrity with minimal degradation, showcasing their resilience to incomplete data [28].

As depicted in Figure 3, the learned causal structure helps security teams trace the origin of a threat signal such as a sequence from stressor events through policy violation and into file access anomalies providing not just predictions, but narrative-driven explanations of risk events [29]. This interpretability ensures PGNs are not only accurate but also operationally trustworthy in high-stakes environments.



Figure 3: Learned Bayesian Network Tracing Causal Pathways from Organizational Stressors to Suspicious File Access with Rounded Posterior Probabilities

This flowchart visualizes a probabilistic graphical model (PGN) illustrating how stressor events lead to policy violations and ultimately to anomalous file access. The edges are annotated with rounded conditional probabilities, enabling interpretable, narrative-driven risk inference.

1. Probability of Transition or Causality

Each decimal shows the likelihood that one event (e.g., stressor) leads to another (e.g., policy violation). For example:

• 0.72 from Stressor Detected → Policy Violation means there's a 72% chance this transition occurs based on historical data or trained model inference.

2. Posterior Probabilities

Some nodes may display a decimal such as 0.58, which could represent:

• The posterior probability that the event (e.g., suspicious file access) occurred given all prior events in the chain.

3. Risk Scores or Trust Scores

In dynamic Bayesian networks or causal inference graphs, these values can also signal:

• Cumulative risk scores or trust degradation metrics calculated in real-time for each user or system state.

6. SIMULATION, CASE STUDY, AND VISUALIZATION

6.1 Organizational Simulation Scenario

To demonstrate the practical application of probabilistic graphical networks (PGNs) for insider threat detection, we simulated an organizational scenario involving a sequence of escalating insider behaviors. The simulation emulated a midsize enterprise environment with 100 employees, integrating event logs that mimic real-world digital and contextual activities.

In the defined scenario, an employee named "User A" received a formal workplace reprimand following a missed project deadline. Subsequently, the simulation introduced behavioral changes over several days, including off-hour logins, repeated access to sensitive customer records, email forwarding to external addresses, and insertion of a personal USB device [23].

These behavioral signals were mapped to causal nodes within the Bayesian network, including "emotional stress," "access anomaly," and "exfiltration intent." As the user's activity deviated from their baseline, the model dynamically updated posterior threat probabilities using belief propagation algorithms. The likelihood of an insider incident increased from 0.07 to 0.81 over four days.

Importantly, the PGN model accounted for latent intent transitions, inferring escalation from normal behavior to high-risk actions based on a combination of digital telemetry and contextual triggers. This progression allowed the model to flag User A's behavior prior to data exfiltration, generating an early warning signal for SOC analysts.

As shown in Figure 4, this simulation provided time-series visualizations of insider likelihood across the user population. Such insights offer operational value, enabling SOCs to prioritize cases that show gradual but converging risk indicators. The scenario validates PGNs as a tool for continuous, causally grounded threat surveillance that adapts to dynamic human behavior [24].

6.2 Case Study: Financial Sector Incident

To assess PGNs in a real-world context, we applied the model to an anonymized dataset representing user activity within a synthetic financial services firm. The dataset included structured logs for login times, transaction approvals, policy exceptions, system alerts, and HR records across 250 employees over three months [25].

Subject Matter Experts (SMEs) from fraud prevention and information security teams contributed to defining expected access roles and escalation paths, which were encoded as domain-informed constraints during the graph learning phase. For instance, a compliance officer accessing trade settlement systems was flagged as a prohibited edge and removed from candidate graph structures [26].

During the study period, a user designated "Employee X" displayed escalating anomalies: early-morning login attempts, unusually high transaction volumes outside job scope, and email exchanges with blacklisted domains. The Bayesian network model identified conditional dependencies between "policy violation," "access role deviation," and "anomalous peer comparison," predicting a 76% likelihood of insider misconduct by week three.

The PGN's predictions were validated post hoc against ground truth labels, which confirmed that Employee X was involved in a policy breach incident that triggered disciplinary proceedings. As depicted in Figure 4, the model provided lead-time visibility into threat evolution, outperforming static rules that failed to recognize the cross-domain behavior pattern.

This case study demonstrates the value of combining PGNs with domain expertise for semantically rich and accurate insider detection, particularly in regulated industries where the cost of delayed response is substantial [27].

6.3 Threat Progression and Early-Warning Visuals

A defining advantage of PGNs is their ability to continuously update posterior probabilities as new behavioral evidence accumulates. This functionality enables visualizations that track insider threat progression across employees and departments over time, facilitating real-time triage and proactive response.

In our simulations, we generated a time-series heatmap (see Figure 4) that visualizes insider threat likelihood scores across organizational units. Each row represents an employee, and each column corresponds to a daily posterior update. Employees with rapidly increasing scores were visually highlighted, allowing analysts to detect emerging threats at a glance [28].

For example, in the simulated enterprise, "User D" exhibited minor anomalies early in the timeline such as non-critical command-line usage. However, subsequent events, including abnormal database queries and unusual file access sequences, caused a cumulative rise in their inferred threat score. By day six, their likelihood score exceeded 0.65, surpassing the predefined risk threshold.

The system generated an alert, automatically annotating contributing nodes in the PGN that influenced the update such as "policy deviation," "temporal access anomaly," and "hidden intent." This traceability helped analysts understand why the model reached a given risk assessment, improving trust and enabling targeted interviews or access revocation [29].

In comparative analysis, PGNs showed faster time-to-detection compared to traditional statistical baselines. As shown in Table 3, the model achieved a detection lead time of 2.8 days, reduced the false positive rate by 19%, and sustained an average posterior update rate of 0.14 per day.

These results underscore the practical advantage of PGNs: rather than reacting to violations, organizations gain a preview of behavioral trajectories, enabling preemptive intervention and reducing both harm and remediation costs [30].



Figure 4: Time-Series Heatmap of Insider Threat Likelihood Scores Across Employees

Figure 4: Time-series heatmap visualizing the evolution of insider threat likelihood scores across employees. Rows represent individual users; columns show daily posterior updates. Redder shades indicate increasing threat probability, highlighting emerging risk clusters by department.

Metric	Value	Description		
Detection Lead Time	2.8 days	Average time PGN flagged a threat before final incident		
False Positive Rate	11.3%	Proportion of flagged users with no actual threat occurrence		
Posterior Update Rate	0.14/day	Average daily change in inferred threat score per user		
Maximum Posterior Score	0.92	Peak likelihood recorded during simulated escalation events		
Prediction Accuracy (F1)	0.87	Balanced score combining precision and recall		

7. INTEGRATION WITH ORGANIZATIONAL TRUST AND CYBER-RESILIENCE STRATEGIES

7.1 Dynamic Trust Scoring Based on Causal Risk

Probabilistic Graphical Networks (PGNs) enable dynamic updates to employee trust scores by leveraging causal relationships inferred from behavior over time. Rather than static background checks or inflexible policy rules, organizations can compute trust as a probabilistic function of contextualized digital behaviors. These scores are derived from posterior probabilities of insider activity, updated in real time as evidence accumulates [27].

Trust scoring integrates directly into Identity and Access Management (IAM) systems. For instance, if an employee's inferred risk score crosses a defined threshold, the IAM system can automatically restrict lateral movement across networks or trigger step-up authentication for sensitive operations. The trust score can also inform access revocation policies during periods of elevated behavioral anomalies [28].

PGNs allow these scores to reflect causal risk rather than superficial correlation. For example, repeated out-of-hours logins may not alone signify risk, but when causally linked to stress indicators and policy deviations, they significantly elevate trust concerns. As shown in Figure 5, this trust architecture includes a feedback loop between PGNs, IAM controls, and HR systems to ensure proportional and explainable responses [29].





Figure 5: Trust Architecture Stack Integrating Causal Modeling, IAM, and HR Oversight

This layered framework illustrates the integration of Probabilistic Graphical Networks (PGNs) with Identity and Access Management (IAM) systems and Human Resources (HR) oversight protocols. The feedback loop ensures that insider threat signals from PGNs inform IAM decisions and are ethically moderated through HR governance. This structure supports explainable, proportional responses while preserving employee trust and organizational transparency.

Furthermore, the trust score is temporally dynamic improving as behavior normalizes or risk nodes deactivate. This adaptive property helps prevent permanent penalization and allows recovery of user privileges over time, aligning security with operational fairness. This dynamic scoring framework also supports behavioral baselining across departments, flagging individuals whose trust deviation diverges from typical peer group patterns [30].

By shifting trust management from rule-based to inference-based, organizations improve their capacity to detect, isolate, and respond to threats proactively, without undermining user productivity or institutional morale.

7.2 Enhancing Organizational Resilience

Beyond threat detection, PGNs offer strategic value in enhancing organizational resilience. These models allow security teams to anticipate not only potential incidents but also the pathways through which disruption spreads informing recovery strategies and resilience planning [31].

For instance, if a PGN identifies causal links between job dissatisfaction, unauthorized access, and potential sabotage, HR and IT departments can jointly simulate policy interventions, such as workload redistribution or role reassignment. This what-if simulation capability allows stakeholders to test how changing a single node like team structure or system permissions affects downstream risk across the organizational graph [32].

This resilience modeling supports prioritization of critical assets and human roles. Employees or departments with high centrality in the PGN meaning their actions affect multiple risk pathways can be proactively audited or supported. In one simulated case, the reassignment of a highly stressed employee to a less sensitive project node reduced network-wide insider threat probability by 18%, as observed in posterior network recalibration [33].

Moreover, PGNs enable rapid post-incident forensic tracing. When an event occurs, analysts can walk back through the network's inferred pathways to identify early causal triggers. This capability shortens incident response cycles and guides the redesign of vulnerable workflows or access policies.

As visualized in Figure 5, these capabilities are integrated into a trust architecture stack, aligning real-time risk monitoring with HR inputs and resilience planning. When threats are seen not as isolated events but as part of evolving behavior chains, organizations can intervene earlier and recover faster, transforming threat modeling into a broader resilience asset [34].

7.3 HR and Ethical Oversight

While PGNs offer potent tools for insider threat mitigation, their application raises ethical and human resources (HR) considerations. Surveillance systems that profile users must include clear oversight protocols to ensure fairness and avoid misuse.

A central principle is the establishment of human-in-the-loop protocols, where flagged threat scores do not immediately trigger punitive actions without HR validation or multi-party review. This helps preserve due process and prevents overreliance on algorithmic judgments [35].

Moreover, privacy-preserving boundaries must be encoded into the modeling framework. PGNs should be restricted from incorporating sensitive personal data such as health or political affiliations unless explicitly permitted and anonymized. Data minimization and legal compliance under frameworks like GDPR must guide both dataset design and model deployment [36].

Bias audits should also be periodically conducted to ensure trust scores do not disproportionately target specific groups or job roles. By integrating ethical oversight into Figure 5's architectural stack, organizations can operationalize AI-based detection without eroding employee trust or workplace cohesion. Ultimately, aligning PGN outputs with transparent and auditable HR governance ensures that security and ethics evolve in parallel, reinforcing both protection and accountability [37].

8. TECHNICAL CHALLENGES AND MITIGATION STRATEGIES

8.1 Data Sparsity and Imbalanced Events

One of the principal challenges in insider threat detection using probabilistic graphical networks (PGNs) is the extreme class imbalance between benign and malicious behaviors. In many enterprise logs, less than 0.5% of events correlate with insider misuse, making it difficult for models to generalize beyond rare examples without overfitting benign behavior [32].

To mitigate this, multiple sampling techniques are employed. Random oversampling of malicious events, while simple, risks duplication artifacts and may inflate false positives. More effective is SMOTE (Synthetic Minority Oversampling Technique), which synthesizes new samples by interpolating minority class examples in feature space [33].

In addition, synthetic data injection has gained traction. Using domain-specific behavioral rules, simulated insider actions such as timed file exfiltration or email misuse are injected into anonymized logs to balance the distribution of causally significant nodes. These synthetic pathways help the PGN learn causal dependencies between hidden stressors and digital outputs without requiring large amounts of real-world breach data [34].

Temporal balancing is also essential. Insider threats often manifest as slow behavioral drift, so sampling strategies must preserve sequence order. Techniques like windowed bootstrapping or causal sequence injection maintain event realism while expanding rare scenarios.

When integrated into the PGN training pipeline, these techniques improve posterior accuracy and resilience against bias. Notably, their effectiveness is amplified when combined with domain-informed structural constraints, which prevent spurious causal edges between over-represented noise events and risk indicators [35]. Thus, intelligent sampling and synthetic augmentation enable PGNs to overcome data sparsity without sacrificing interpretability or robustness.

8.2 Scalability and Model Drift

As organizations scale, PGNs face challenges in handling data from thousands of users and millions of interactions. Additionally, insider behaviors evolve what constitutes risk today may be benign tomorrow. These factors introduce the risk of model drift, where predictive accuracy deteriorates over time [36].

To address this, many implementations rely on incremental learning strategies. Rather than retraining the entire PGN from scratch, models are updated using Bayesian online learning, where parameter posteriors are refreshed as new evidence arrives. This ensures that trust scores and risk pathways evolve with organizational dynamics [35].

Batch updates, performed weekly or monthly, offer stability but may lag behind fast-developing threats. Conversely, real-time updates using streaming logs allow immediate risk recalibration but can introduce volatility if not dampened by priors or temporal smoothing [37].

Scalable implementations also involve distributed learning, where subnetworks are trained locally within departments and later merged using structure consensus protocols. This modular approach enables faster updates, reduces central compute load, and preserves interpretability at the subgraph level.

In practice, scalability requires not just computational optimization, but also governance to monitor drift thresholds and update triggers, ensuring the PGN adapts to new risk patterns while avoiding overfitting or instability.

8.3 Interpretability vs. Complexity

A final operational consideration is the trade-off between model complexity and interpretability, especially when PGNs are used by non-technical stakeholders such as HR, compliance teams, and executives. High-dimensional networks with dozens of causal variables may capture nuanced behaviors but often become difficult to explain [38].

Interpretability is critical when PGNs inform access revocation, HR action, or compliance reporting. A model that flags an employee must also indicate why which pathways were activated, what causal links were observed, and how confident the inference was. As seen in Figure 5, effective architectures embed explanation layers that present simplified subgraphs with annotated node transitions.

To manage this trade-off, one approach is hierarchical modeling, where macro-behaviors (e.g., "anomalous access") are inferred first, followed by detailed subgraphs when needed. This keeps the primary alert layer lightweight and interpretable, while deeper investigation tools remain available for forensic use [39].

Another method involves saliency analysis ranking nodes or edges based on their contribution to threat probability. This helps analysts prioritize attention and understand model rationale without interpreting the full graph. Rule extraction techniques, derived from learned graph structures, can also translate probabilistic inferences into decision-tree–like explanations for compliance audits.

Ultimately, the goal is not merely predictive accuracy but actionable transparency. PGNs can retain their causal modeling power while offering digestible outputs for human decision-makers through modularization, visualization, and strategic simplification. This ensures alignment with real-world trust, fairness, and usability goals [40].

9. FUTURE DIRECTIONS AND CROSS-DOMAIN APPLICATIONS

9.1 Integration with Federated Learning Systems

Probabilistic Graphical Networks (PGNs) can be seamlessly integrated into federated learning (FL) frameworks to enable distributed causal inference across multi-national or multi-organizational infrastructures. This is particularly valuable for global enterprises seeking to identify behavioral threats without centralizing sensitive employee or contractor data [35].

Under this paradigm, each regional node or division trains local PGNs using anonymized and differentially private logs, retaining jurisdictional data sovereignty. Once trained, model parameters not raw data are shared and aggregated at a central orchestrator using secure protocols such as SMPC or homomorphic encryption [36].

This decentralized learning architecture enables unified threat pattern recognition while respecting privacy and compliance frameworks like GDPR or HIPAA. For example, a Europe-based HR department may detect a causal structure linking user stressors to anomalous data movements, which, once aggregated, enhances global predictive accuracy without disclosing identifiable logs [32].

Such integration enhances scalability while allowing for cross-border insider threat detection, creating a globally distributed, privacy-preserving causal intelligence mesh. It also supports adaptive trust scoring in multinational firms, where behavioral context varies but risk indicators remain causally aligned [37].

9.2 PGNs for Supply Chain Insider Risk

Modern supply chains rely on complex networks of external vendors, contractors, and third-party service providers introducing substantial insider risk beyond organizational boundaries. Traditional security models often fail to account for behavioral anomalies arising from external partners embedded in logistics, IT, or manufacturing processes [38].

PGNs offer a compelling framework for zero-trust modeling, where every actor internal or external is evaluated based on observed causal behavior rather than presumed affiliation. By incorporating access frequency, credential movement, and interface misuse, PGNs can flag subcontractors exhibiting behavior that deviates from established trust baselines.

For example, a logistics vendor repeatedly accessing time-sensitive production schedules outside expected hours could activate nodes linked to financial motive and schedule tampering risk. This event chain becomes even more actionable when aligned with contextual signals like recent contract disputes or project overruns [39].

As shown in Table 1, PGNs provide higher-resolution detection granularity compared to conventional access-control logs. Moreover, they support contractor-specific policy refinement, where access rights and risk thresholds dynamically update as causal insights evolve. These models thus become essential tools for supply chain assurance, enabling real-time monitoring of third-party behavior within a federated risk management framework [40].

9.3 Cyber-Physical Systems and Critical Infrastructure

Causal modeling is particularly effective in Cyber-Physical Systems (CPS) where logical behaviors and physical outcomes are tightly coupled. In critical infrastructure environments like energy grids, water treatment plants, or industrial automation, insider attacks may manifest as minor command deviations that cascade into operational sabotage [41].

PGNs enable the discovery of behavioral precursors to such events. For instance, a technician bypassing safety checks before adjusting turbine control parameters may trigger causal alerts linking the behavior to stress indicators, privilege escalation patterns, or remote login anomalies. Temporal Bayesian models can track how such behaviors unfold over time critical in environments where milliseconds matter [42].

These insights also support multi-layered defense strategies, where behavioral logs (keystrokes, badge access) are fused with sensor data and operational metrics (e.g., pressure differentials, PLC commands). When integrated with SCADA threat monitoring, PGNs act as semantic interpreters, helping SOC teams understand not just what failed, but why [43].

As illustrated in Figure 4, real-time causal overlays on control systems allow for early warnings before catastrophic failure or compromise. By embedding PGNs within CPS security architecture, critical infrastructure gains interpretable, predictive, and actionable defense against human-driven sabotage or insider lapses [44].

9.4 AI-Augmented Behavioral Defense Teams

Future SOCs (Security Operations Centers) are increasingly envisioned as hybrid human-AI teams, where analysts collaborate with intelligent assistants that provide real-time behavioral context. PGNs serve as foundational inference engines for such teams, offering causal narratives rather than isolated alerts [45].

When a potential insider event occurs say, unauthorized document access at midnight the PGN surfaces not just the event but its causal roots: workload complaints, reduced collaboration, anomalous browsing history. These contextual paths guide SOC analysts in distinguishing between false alarms and legitimate threats [46].

Moreover, the PGN's probabilistic scoring helps prioritize alerts. Rather than overwhelming analysts with raw log streams, the system ranks employees or nodes based on cumulative risk trajectories. Visualizations, such as those in Figure 5, show evolving threat paths and confidence levels, enabling more focused investigation [47].

SOC workflows also benefit from feedback loops. Analysts can confirm or dismiss inferred threats, improving the model's future predictions. Over time, the PGN learns organization-specific risk dynamics becoming a living behavioral map that strengthens both detection and trust [48].

This symbiosis between humans and AI, anchored by causal reasoning, enables faster, fairer, and more precise cyber defense, transforming SOCs into adaptive behavioral intelligence hubs.

10. CONCLUSION

10.1 The Strategic Advantage of Causal Modeling with PGNs for Insider Threat Detection

In an increasingly complex cybersecurity landscape, organizations face the dual challenge of identifying nuanced insider threats while maintaining operational trust and ethical oversight. Probabilistic Graphical Networks (PGNs) offer a transformative capability to meet these demands, moving beyond traditional anomaly detection or rule-based methods to establish a foundation rooted in causality and interpretability.

Unlike conventional models that rely heavily on surface-level pattern recognition, PGNs uncover the underlying structure of insider behaviors. By modeling relationships between variables such as stress levels, access patterns, digital anomalies, and workplace context, PGNs provide a transparent map of how threats emerge not just when or where. This causal perspective enables a more anticipatory approach, flagging risky behavior chains long before they escalate into damaging breaches. Importantly, it accounts for both malicious insiders and unintentional actors, such as employees under duress who might inadvertently violate security protocols.

A core strength of PGNs lies in their ability to contextualize digital behaviors within organizational dynamics. These models do not operate in isolation but instead integrate with access logs, HR metadata, workflow systems, and identity frameworks to infer intent and motivation. The ability to model latent constructs such as stress, fatigue, or grievance based on observable signals sets PGNs apart as intelligent systems that mimic human reasoning while remaining statistically grounded.

This dual capability of predictive accuracy and interpretability makes PGNs especially valuable for decision-makers. Unlike deep learning models, which often function as "black boxes," PGNs can articulate why a certain behavior was flagged, which causal pathways were activated, and what the probable outcomes might be. This transparency builds institutional trust, enabling managers, HR teams, and compliance officers to act confidently and fairly on model outputs. It also supports ethical implementation, allowing for justifiable scrutiny and auditability in sensitive use cases such as access revocation or internal investigations.

Crucially, PGNs offer a blueprint for human-machine hybrid trust architectures. In this paradigm, algorithms serve not as authoritarian detectors but as advisory partners. They augment security operations by surfacing high-risk behaviors and guiding investigation paths, while humans provide judgment, nuance, and empathy. These hybrid workflows preserve organizational fairness and resilience, avoiding overreliance on either automation or subjective decision-making.

Moreover, PGNs scale naturally into federated infrastructures, supporting multinational operations without compromising data sovereignty. Their ability to operate within federated learning environments allows behavioral models to be trained on distributed logs while maintaining compliance with privacy regulations. This ensures organizations can coordinate insider threat detection across global branches, vendors, and supply chains in a privacy-preserving, policy-aligned manner.

Looking ahead, adopting causal-aware frameworks is not merely a technical upgrade it is a strategic imperative. As digital ecosystems expand and insider threats grow more sophisticated, reactive security postures will no longer suffice. Organizations must move toward proactive, explainable, and context-rich modeling approaches that align security with organizational psychology, ethics, and operational integrity.

PGNs represent a step toward this vision: an interdisciplinary, resilient, and human-centric defense architecture. Their capacity to uncover why threats happen rather than just when marks a profound shift in cyber defense philosophy. By embedding PGNs into trust scoring systems, incident response protocols, and access governance workflows, organizations can evolve from static monitoring to adaptive, causal intelligence systems.

In conclusion, the path to sustainable cyber-resilience lies in fusing the predictive power of PGNs with the discernment of human analysts. Causal modeling empowers organizations not just to detect threats but to understand them, manage them, and prevent them with clarity, confidence, and accountability.

REFERENCE

- Ademilua DA, Areghan E. AI-Driven Cloud Security Frameworks: Techniques, Challenges, and Lessons from Case Studies. Communication In Physical Sciences. 2022 Dec 30;8(4):674-88.
- Joshi A, Chauhan N, Thakur G, Kumar V, Singh Y. Fortifying Tomorrow: Overcoming Challenges in Cloud AI Ecosystems. InDeep Learning Innovations for Securing Critical Infrastructures 2025 (pp. 507-526). IGI Global Scientific Publishing.
- Jamiu OA, Chukwunweike J. DEVELOPING SCALABLE DATA PIPELINES FOR REAL-TIME ANOMALY DETECTION IN INDUSTRIAL IOT SENSOR NETWORKS. International Journal Of Engineering Technology Research & Management (IJETRM). 2023Dec21;07(12):497–513.
- Talla RR, Manikyala A, Nizamuddin M, Kommineni HP, Kothapalli S, Kamisetty A. Intelligent Threat Identification System: Implementing Multi-Layer Security Networks in Cloud Environments. NEXG AI Review of America. 2021;2(1):17-31.
- 5. Emehin O, Akanbi I, Emeteveke I, Adeyeye OJ. Enhancing cybersecurity with safe and reliable AI: mitigating threats while ensuring privacy protection. International Journal of Computer Applications Technology and Research, doi. 2024;10.
- Tabrizchi H, Aghasi A. Cyber Security Intelligent Systems Based on Federated Learning. InFederated Cyber Intelligence: Federated Learning for Cybersecurity 2025 Apr 24 (pp. 75-100). Cham: Springer Nature Switzerland.
- Whig P, Aggarwal A, Ganeshan V, Modhugu VR, Bhatia AB. AI for Secure and Resilient Cyber-Physical Systems. InArtificial Intelligence Solutions for Cyber-Physical Systems 2024 (pp. 40-63). Auerbach Publications.
- 8. Hussain A. AI and Machine Learning in Action: Revolutionizing Enterprise Data Security and Cloud Infrastructure Protection. Journal of Cloud Computing and AI Integration. 2024 Nov.
- Gangapatnam K. Proactive Security with AI: Revolutionizing Cloud Infrastructure Protection. Journal of Computer Science and Technology Studies. 2025 May 3;7(3):277-84.
- Odumbo OR, Nimma SZ. Leveraging artificial intelligence to maximize efficiency in supply chain process optimization. Int J Res Publ Rev. 2025;6(01):[pages not specified]. doi: <u>https://doi.org/10.55248/gengpi.6.0125.0508</u>.
- Unanah Onyekachukwu Victor, Yunana Agwanje Parah. Clinic-owned medically integrated dispensaries in the United States; regulatory pathways, digital workflow integration, and cost-benefit impact on patient adherence (2024). *International Journal of Engineering Technology Research & Management (IJETRM)*. Available from: https://doi.org/10.5281/zenodo.15813306
- Sundaramurthy SK, Ravichandran N, Inaganti AC, Muppalaneni R. AI-Driven Threat Detection: Leveraging Machine Learning for Real-Time Cybersecurity in Cloud Environments. Artificial Intelligence and Machine Learning Review. 2025 Jan 15;6(1):23-43.
- Anandharaj N. AI-powered cloud security: A study on the integration of artificial intelligence and machine learning for improved threat detection and prevention. J. Recent Trends Comput. Sci. Eng.(JRTCSE). 2024 Jul 25;12:21-30.
- 14. Lad S. Cybersecurity trends: Integrating ai to combat emerging threats in the cloud era. Integrated Journal of Science and Technology. 2024 Aug 12;1(3).

- 15. Singh K, Saxena R, Kumar B. AI Security: Cyber Threats and Threat-Informed Defense. In2024 8th Cyber Security in Networking Conference (CSNet) 2024 Dec 4 (pp. 305-312). IEEE.
- Berna IE, Vijay K, Gnanavel S, Jeyalakshmi J. Impact of Artificial Intelligence and Machine Learning in Cloud Security. InImproving Security, Privacy, and Trust in Cloud Computing 2024 (pp. 34-58). IGI Global Scientific Publishing.
- Anwar N, Rahaman R, Widodo AM, Sekti BA, Erzed N, Tangkudung RR, Muchlis M, Yuhefizar Y, Budhisantosa N. Sustainable Cybersecurity in the AI Era: Innovations in Green Computing and Security. InSustainable Information Security in the Age of AI and Green Computing 2025 (pp. 603-620). IGI Global Scientific Publishing.
- Sharma DP, Habibi Lashkari A, Firoozjaei MD, Mahdavifar S, Xiong P. AI for Cloud Security. InUnderstanding AI in Cybersecurity and Secure AI: Challenges, Strategies and Trends 2025 May 27 (pp. 95-111). Cham: Springer Nature Switzerland.
- Ali A, Razzaque A, Munir U, Shahid H, Khattak FW, Rajpoot Z, Kamran M, Farid Z. AI-Driven Approaches to Cybersecurity: The Impact of Machine and Deep Learning. In2024 2nd International Conference on Cyber Resilience (ICCR) 2024 Feb 26 (pp. 1-5). IEEE.
- Arif A, Khan MI, Khan AR, Anjum N, Arif H. AI-Driven Cybersecurity Predictions: Safeguarding California's Digital Landscape. International Journal of Innovative Research in Computer Science and Technology. 2025;13:74-8.
- Shaffi SM, Vengathattil S, Sidhick JN, Vijayan R. AI-Driven Security in Cloud Computing: Enhancing Threat Detection, Automated Response, and Cyber Resilience. arXiv preprint arXiv:2505.03945. 2025 May 6.
- Celeste R, Michael S. Next-Gen Network Security: Harnessing AI, Zero Trust, and Cloud-Native Solutions to Combat Evolving Cyber Threats. International Journal of Trend in Scientific Research and Development. 2021;5(6):2056-69.
- Huma Z, Muzaffar J. Hybrid AI Models for Enhanced Network Security: Combining Rule-Based and Learning-Based Approaches. Global Perspectives on Multidisciplinary Research. 2024 Sep 30;5(3):52-63.
- 24. Zhang X, Wang P, Jia H, Huang Z, Zhao R. AI-Powered Cybersecurity: Enhancing Threat Detection and Defense in the Digital Age. In2024 IEEE 7th International Conference on Electronic Information and Communication Technology (ICEICT) 2024 Jul 31 (pp. 1026-1031). IEEE.
- 25. Raza H. Proactive cyber defense with AI: Enhancing risk assessment and threat detection in cybersecurity ecosystems. Journal Name Missing. 2021 Jul 11.
- Sobur A, Hossain A. MACHINE LEARNING IN CYBERSECURITY: HARNESSING AI FOR DEFENSE. Available at SSRN 4878189. 2024 Mar 3.
- 27. Abubakar M, Chitraju Gopal Varma S, Volikatla H, Likki H, MS H. Leveraging AI and Machine Learning for Enhanced Cloud Security and Performance. Hemanth and Likki, Hemanth and gp, hemanth and S, Hemanth and MS, Hemanth, Leveraging AI and Machine Learning for Enhanced Cloud Security and Performance (May 14, 2020). 2020 May 14.
- Laura M, James A. Cloud Security Mastery: Integrating Firewalls and AI-Powered Defenses for Enterprise Protection. International Journal of Trend in Scientific Research and Development. 2019;3(3):2000-7.
- 29. Pemmasani PK. AI in National Security: Leveraging Machine Learning for Threat Intelligence and Response. The Computertech. 2023 Jan 19:1-0.
- Shahana A, Hasan R, Farabi SF, Akter J, Mahmud MA, Johora FT, Suzer G. AI-driven cybersecurity: Balancing advancements and safeguards. Journal of Computer Science and Technology Studies. 2024 May 10;6(2):76-85.
- Nina P, Ethan K. AI-driven threat detection: Enhancing cloud security with cutting-edge technologies. International Journal of Trend in Scientific Research and Development. 2019;4(1):1362-74.
- 32. Tarafdar R. AI-Powered Cybersecurity Threat Detection in Cloud Environments. International Journal of Cybersecurity and Digital Forensics. 2022.
- Ofili BT, Obasuyi OT, Osaruwenese E. Threat intelligence and predictive analytics in USA cloud security: mitigating AI-driven cyber threats. Int J Eng Technol Res Manag. 2024 Nov;8(11):631.
- 34. Rehan H. AI-driven cloud security: The future of safeguarding sensitive data in the digital age. Journal of Artificial Intelligence General science (JAIGS) ISSN. 2024 Jan: 3006-4023.
- 35. Sharma M, Raymond D, Weerarathna I, Kumar P, Chadar AR. Cloud Security and Artificial Intelligence. InHandbook of AI-Driven Threat Detection and Prevention (pp. 144-164). CRC Press.
- 36. Omar M, Zangana HM, editors. Redefining Security With Cyber AI. IGI Global; 2024 Jul 17.
- Aarav M, Layla R. Cybersecurity in the cloud era: Integrating AI, firewalls, and engineering for robust protection. International Journal of Trend in Scientific Research and Development. 2019;3(4):1892-9.

- Machhindra PA, Vijay BN, Mahendra BS, Rahul CA, Anil PA, Sunil PR. Enhancing cyber security through machine learning: A comprehensive analysis. In2023 4th International Conference on Computation, Automation and Knowledge Management (ICCAKM) 2023 Dec 12 (pp. 1-6). IEEE.
- Bhambri P, Bajdor P. AI-Powered Predictive Analysis for Proactive Cyber Defense. InHandbook of AI-Driven Threat Detection and Prevention (pp. 308-321). CRC Press.
- 40. Sharma A, Bhatia R, Sharma D, Kalra A. Exploring Al's Prowess in Advancing Cybersecurity. InSmart Systems: Engineering and Managing Information for Future Success: Navigating the Landscape of Intelligent Technologies 2025 Feb 25 (pp. 77-98). Cham: Springer Nature Switzerland.
- 41. Bhambri P. Understanding AI and Machine Learning in Security. InHandbook of AI-Driven Threat Detection and Prevention 2025 (pp. 1-17). CRC Press.
- 42. Bolanle O, Bamigboye K. AI-Powered Cloud Security: Leveraging Advanced Threat Detection for Maximum Protection. International Journal of Trend in Scientific Research and Development. 2019;3(2):1407-12.
- 43. Lekkala S, Gurijala P. Leveraging AI and Machine Learning for Cyber Defense. InSecurity and Privacy for Modern Networks: Strategies and Insights for Safeguarding Digital Infrastructures 2024 Oct 8 (pp. 167-179). Berkeley, CA: Apress.
- 44. Manojkumar R, Vaidya SP, Jena PK, Ashok S, Dasari S. Machine Learning and Federated Learning in Industrial Cybersecurity. InAI-Enhanced Cybersecurity for Industrial Automation 2025 (pp. 407-438). IGI Global Scientific Publishing.
- 45. Muppalaneni R, Inaganti AC, Ravichandran N. AI-Driven Threat Intelligence: Enhancing Cyber Defense with Machine Learning. Journal of Computing Innovations and Applications. 2024 Jan 12;2(1):1-1.
- 46. Emmanni PS. Federated learning for cybersecurity in edge and cloud computing. International Journal of Computing and Engineering. 2021;2(1):14-25.
- 47. Rahman MK, Dalim HM, Hossain MS. AI-Powered solutions for enhancing national cybersecurity: predictive analytics and threat mitigation. International Journal of Machine Learning Research in Cybersecurity and Artificial Intelligence. 2023;14(1):1036-69.
- Singh SK, Bhambu P, Sandhu A, Kumar A, Sharma D, Pandey A. Achieving Cloud Security Solutions based on Machine Learning and Past Information. In2024 International Conference on Augmented Reality, Intelligent Systems, and Industrial Automation (ARIIA) 2024 Dec 20 (pp. 1-6). IEEE.